

E0917/95-319

APPARATUS AND METHOD FOR IDENTIFYING DATA  
PACKET AT WIRE RATE ON A NETWORK  
SWITCH PORT

BACKGROUND OF THE INVENTION

FIELD OF THE INVENTION

The present invention relates to layer 2 and layer 3 switching of data packets in a non-blocking network switch configured for switching data packets between subnetworks.

BACKGROUND ART

5 Local area networks use a network cable or other media to link stations on the network. Each local area network architecture uses a media access control (MAC) enabling network interface devices at each network node to access the network medium.

The Ethernet protocol IEEE 802.3 has evolved to specify a half-duplex media access mechanism and a full-duplex media access mechanism for transmission of data packets. The full-  
10 duplex media access mechanism provides a two-way, point-to-point communication link between two network elements, for example between a network node and a switched hub.

Switched local area networks are encountering increasing demands for higher speed connectivity, more flexible switching performance, and the ability to accommodate more complex network architectures. For example, commonly-assigned U.S. Patent No. 5,953,335 discloses a  
15 network switch configured for switching layer 2 type Ethernet (IEEE 802.3) data packets between different network nodes; a received data packet may include a VLAN (virtual LAN) tagged frame according to IEEE 802.1q protocol that specifies another subnetwork (via a router) or a prescribed group of stations. Since the switching occurs at the layer 2 level, a router is typically necessary to transfer the data packet between subnetworks.

20 Efforts to enhance the switching performance of a network switch to include layer 3 (e.g., Internet protocol) processing may suffer serious drawbacks, as current layer 2 switches preferably are configured for operating in a non-blocking mode, where data packets can be output from the switch at the same rate that the data packets are received. Newer designs are needed to ensure that higher speed switches can provide both layer 2 switching and layer 3 switching capabilities for faster speed  
25 networks such as 100 Mbps or gigabit networks.

004780-5702E960

However, such design requirements risk loss of the non-blocking features of the network switch, as it becomes increasingly difficult for the switching fabric of a network switch to be able to perform layer 3 processing at the wire rates (i.e., the network data rate).

Conventional routers and high-end switches use a centralized software based CPU intensive scheme to examine data packets for classifications and to determine the class or priority for switching/routing the packet. Figure 1 shows the time sequence of a data packet as it arrives at a conventional network switch port (T1) until it leaves the network switch on an output port (T10). Since the packet must be sent to a processor (at T3) and then classified by the processor (at T5-T6), the latency (first bit in to first bit out) is high and can be expressed as  $S1+S2+S3+S4$ . The latency is high since processing of the packet does not even begin until it is fully received (after T2).

### SUMMARY OF THE INVENTION

There is a need for an arrangement that enables a network switch to provide layer 2 switching and layer 3 switching capabilities for 100 Mbps and gigabit links without blocking of the data packets.

There is also a need for an arrangement that enables a network switch to provide layer 2 switching and layer 3 switching capabilities with minimal buffering within the network switch that may otherwise affect latency of switched data packets.

There is also a need for an arrangement that enables a network switch to be easily programmable to distinguish between different types of layer 3 data packets so that quality of service (QoS) can be achieved.

There is also a need for an arrangement to enable a network switch port to instantaneously evaluate an incoming data packet and determine a layer 3 or higher protocol, to provide the associated switch fabric with sufficient time to process the incoming data packet according to the detected protocol.

These and other needs are attained by the present invention, where, in a network switch port, a frame tag is generated based on the evaluation of an entire data packet as soon as a last bit of the data packet is received at the network switch port.

One aspect of the present invention provides a method of evaluating an incoming data packet at a network switch port. The method includes storing a plurality of templates configured for identifying respective data formats, each template having at least one min term configured for comparing a corresponding prescribed value to a corresponding selected byte of the incoming data packet. The min terms that correspond to the selected byte are simultaneously compared to the selected byte immediately upon receipt of the selected byte by the network switch port. A comparison

004780" 5707E960

result is generated that identifies the incoming data packet, based on the comparisons of the min terms to the data bytes of the entire data packet received by the network switch port. A frame tag is generated based on the comparison result as soon as a last bit of the data packet is received at the network switch port.

5 Additional advantages and novel features of the invention will be set forth in part in the description which follows and in part will become apparent to those skilled in the art upon examination of the following or may be learned by practice of the invention. The advantages of the present invention may be realized and attained by means of instrumentalities and combinations particularly pointed in the appended claims.

## 10 BRIEF DESCRIPTION OF THE DRAWINGS

Reference is made to the attached drawings, wherein elements having the same reference numeral designations represent like element elements throughout and wherein:

Figure 1 is a diagram showing latency of a data packet defined from the time the packet arrives at a conventional network switch port until the packet leaves the port.

15 Figure 2 is a block diagram of a packet switched network including multiple network switches for switching data packets between respective subnetworks according to an embodiment of the present invention.

Figure 3 is a diagram illustrating a conventional layer 2 Ethernet-type data packet carrying a layer 3 Internet protocol (IP) packet.

20 Figure 4 is a flow diagram illustrating a conventional (prior art) method of evaluating an IP packet.

Figure 5 is a block diagram illustrating the network switch port of the network of Figure 2, including the network switch port filter, according to an embodiment of the present invention.

25 Figures 6A and 6B are diagrams illustrating simultaneous processing of four templates of an equation by the min term generator of Figure 5.

Figure 7 is a diagram illustrating in further detail the simultaneous processing of min terms by the min term generator of Figure 5.

Figure 8 is diagram illustrating a table entry of a min term in the min term memory of Figure 5 according to an embodiment of the present invention.

30 Figure 9 is a diagram showing latency of a data packet of an embodiment of the present invention defined from the time the packet arrives at a network switch port until the packet leaves the port.

004780" 5T02E960  
09637015" 081400

## BEST MODE FOR CARRYING OUT THE INVENTION

Figure 2 is a block diagram illustrating a packet switched network 10, such as an Ethernet (IEEE 802.3) network. The packet switched network includes integrated (i.e., single chip) multiport switches 12 that enable communication of data packets between network stations 14. Each network station 14, for example a client workstation, is typically configured for sending and receiving data packets at 10 Mbps or 100 Mbps according to IEEE 802.3 protocol. Each of the integrated multiport switches 12 are interconnected by gigabit Ethernet links 16, enabling transfer of data packets between subnetworks 18a, 18b, and 18c. Hence, each subnetwork includes a switch 12, and an associated group of network stations 14.

Each switch 12 includes a switch port 20 that includes a media access control (MAC) module 22 and a packet classifier module 24. The MAC module 20 transmits and receives data packets to the associated network stations 14 across 10/100 Mbps physical layer (PHY) transceivers (not shown) according to IEEE 802.3u protocol. Each switch 12 also includes a switch fabric 25 configured for making frame forwarding decisions for received data packets. In particular, the switch fabric 25 is configured for layer 2 switching decisions based on source address, destination address, and VLAN information within the Ethernet (IEEE 802.3) header; the switch fabric 25 is also configured for selective layer 3 switching decisions based on evaluation of an IP data packet within the Ethernet packet.

As shown in Figure 2, each switch 12 has an associated host CPU 26 and a buffer memory 28, for example an SSRAM. The host CPU 26 controls the overall operations of the corresponding switch 12, including programming of the switch fabric 25. The buffer memory 28 is used by the corresponding switch 12 to store data frames while the switch fabric 25 is processing forwarding decisions for the received data packets.

As described above, the switch fabric 25 is configured for performing layer 2 switching decisions and layer 3 switching decisions. The availability of layer 3 switching decisions may be particularly effective if an end station 14 within subnetwork 18a wishes to send an e-mail message to selected network stations in subnetwork 18b, 18c, or both; if only layer 2 switching decisions were available, then the switch fabric 25 of switch 12a would send the e-mail message to switches 12b and 12c without specific destination address information, causing switches 12b and 12c to flood all their ports. Otherwise, the switch fabric 25 of switch 12a would need to send the e-mail message to a router (not shown), which would introduce additional delay. Use of layer 3 switching decisions by the switch fabric 25 enables the switch fabric 25 to make intelligent decisions as far as how to handle a packet, including advanced forwarding decisions, and whether a packet should be considered a high-priority packet for latency-sensitive applications, such as video or voice. Use of layer 3 switching decisions by the switch fabric 25 also enables the host CPU 26 of switch 12a to remotely program another switch, for

004780" 5704960

example switch 12b, by sending a message having an IP address corresponding to the IP address of the switch 12b; the switch 12b, in response to detecting a message addressed to the switch 12b, can forward the message to the corresponding host CPU 26 for programming of the switch 12b.

Figure 3 is a diagram illustrating an Ethernet (IEEE 802.3) packet 30 carrying an IP packet 32 as payload data. Specifically, the Ethernet packet 30 includes a start frame delimiter (SFD) 34, an Ethernet header 36, the IP packet 32, and a cyclic redundancy check (CRC) or frame check sequence (FCS) field 38. Hence, a switch fabric 25 configured for layer 3 switching decisions needs to be able to quickly process the IP packet 32 within the received Ethernet frame 30 to avoid blocking of the frame within the switch.

Figure 4 is a flow diagram illustrating an example of the type of layer 3 processing that might be performed for an incoming data packet. The flow diagram of Figure 3, conventionally implemented in software, would involve checking whether the incoming data packet was a hypertext transport protocol (HTTP) packet in step 50, an SNMP packet in step 52, or a high-priority packet in step 54. The appropriate tag would then be assigned identifying the packet in steps 56, 57, 58, or 60.

The arrangement of Figure 4, however, cannot from a practical standpoint be implemented in hardware in a manner that would provide a non-blocking switch for 100 Mbps or gigabit networks. In particular, the sequential nature of the decision process in Figure 3 would result in undue latency for the incoming data packet.

According to the disclosed embodiment, the packet classifier module 24 of Figure 2 is configured for multiple simultaneous comparisons between the incoming data stream and templates that identify the data format of the incoming data stream. Specifically, users of the host processor 26 will specify policies that define how data packets having certain IP protocols should be handled by the switch fabric 25. These policies are implemented by loading into the switch fabric 25 a set of frame forwarding decisions for each corresponding IP protocol type. Hence, the switch fabric 25 could include one set of frame forwarding instructions for an HTTP packet, another set of frame forwarding instructions for an SNMP packet, and another set of frame forwarding instructions for a high-priority packet (e.g., video, or voice, etc.).

Figure 5 is a block diagram illustrating the packet classifier module 24 according to an embodiment of the present invention. As shown in Figure 5, the network switch port 20 includes a MAC 22, a receive FIFO buffer 27, a header modifier 29, and the packet classifier module 24. The packet classifier module 24, also referred to as a network switch port filter, is configured for identifying (i.e., evaluating) the incoming data packet at the network switch port 20, and supplying to the switch fabric 25 a tag that specifies the action to be performed on the data packet based on type of data packet being received. Specifically, and as shown in Figure 9, the packet classifier module 24 simultaneously

004780" ST04E960

compares the incoming data packet with a plurality of templates configured for identifying respective data formats immediately as the data is received at the switch port 20. The packet classifier module 24, based on the comparison between the incoming data packet and the plurality of templates, identifies an equation to be executed that specifies the tag to be supplied to the switch fabric 25. As shown in Figure 9, since the packet is immediately classified upon reception at the switch port 20, regardless of whether the packet is valid, the latency,  $S1+S2+S3$ , is reduced as compared to the latency described above with reference to Figure 1.

Specifically, the packet classifier module 24 generates a comparison result that identifies the incoming data packet by detecting at least one matched template from a plurality of templates. The packet classifier module 24 then identifies which of the equations includes the matched template, and generates the tag specified by the equation.

Figures 6A and 6B are diagrams illustrating the simultaneous processing of two templates of an equation by the packet classifier module 24. Figure 6A illustrates the logical evaluation by the packet classifier module 24 of the equation:

$$\text{Eq1} = M1 * M2 * M3 * M4 * (M5 + M6 + M7 + M8).$$

Figure 6B illustrates how the equation Eq1 would actually be stored in the min term memory 70. The equation Eq1 includes four templates 62a, 62b, 62c, and 62d: the template 62a includes the min terms M1, M2, M3, M4, and M5; the template 62b includes the min terms M1, M2, M3, M4, and M6; the template 62c includes the min terms M1, M2, M3, M4, and M7; and the template 62d includes the min terms M1, M2, M3, M4, and M8. Each template 62 corresponds to a specific IP data format recognizable based on the header of the IP data packet 32. For example, templates 62a and 62c may be configured for identifying an HTTP packet, and templates 62b and 62d may be configured for identifying an SNMP packet. Specifically, an HTTP packet is identified if it is in IPv4 format, the time to live field in IP is bigger than one, the protocol field in IP header is TCP, header checksum is correct, source TCP port is 80 or destination TCP port is 80. An SNMP packet is identified if it is in IPv4 format, the time to live field in IP is bigger than one, the protocol field in IP header is TCP, header checksum is correct, source TCP port is 25 or destination TCP port is 25.

Hence, the following min terms may be established to represent all the above-described criteria:

- M1= packet is in IPv4 format
- M2= time to live field in IP is bigger than one
- M3= protocol field in IP header is TCP
- M4= header checksum is correct
- M5= source TCP port is 80

M6= destination TCP port is 80

M7= source TCP port is 25

M8= destination TCP port is 25

5 Hence, the templates 62a and 62c identify HTTP packets, and the templates 62b and 62d identify SNMP packets. Thus, equation one (Eq1) specifies that a specific result (e.g., the tag having a specified value) should be output to the switch fabric 25 if either template 62a, 62b, 62c, or 62d are true.

Moreover, the min terms M1...M8 are arranged within the associated templates 62a and/or 62b in a prescribed order that corresponds to the relative position of a data byte in the incoming data stream. As illustrated in Figure 7, the min term M1 is configured for comparison with the first byte (B1) of the IP packet 32, the min term M2 is configured for comparison with a subsequent byte (B2) of the IP packet 32 that follows B1, the min term M3 is configured for comparison with a subsequent byte (B3) that follows B2, etc. Hence, the use of templates 62 having min terms in an order based on the relative position of a data byte in the incoming data stream enables multiple simultaneous comparisons between the incoming data stream and min terms. Hence, an incoming data packet can be compared to multiple templates to determine not only the data format of the incoming data packet, but also what action needs to be performed by the switch fabric 25.

As shown in Figure 5, the packet classifier 24, also referred to as a network switch port filter, includes a min term memory 70 for storing the min term values (e.g., M1, M2, etc.) as illustrated in Figure 8, described below. The packet classifier 24 also includes a frame identifier 72 configured for identifying the type of layer 2 frame being received; in particular, identifying the type of layer 2 frame being received (e.g., Ethernet, IEEE 802 to 3, etc.) enables identification of the start position 64 of the IP packet 32 within the layer 2 packet 30. The packet classifier 24 also includes a min term controller 74, a min term generator 76, an equation core 78, and an evaluation results memory 80. A processor interface module (pi\_mod) 82 is used for transferring the generated min terms from the host CPU 26 into the min term memory 70.

The min term controller 74 is configured for fetching the min terms from the min term memory 70 corresponding to a selected byte of the IP frame 32. The min term controller 74 also includes a location converter configured for specifying the actual byte location (byte\_location) of the start point 64 in response to receiving a frame type (frm\_type) signal from the frame identifier 72 that specifies the type of layer 2 frame. Hence, the min term controller 74, in response to detecting the beginning of the IP packet, fetches all the min terms that are to be compared with the first byte (B1) of the IP packet 32, for example min terms M1, M9, and M14 for equations Eq1, Eq2, and Eq3 in Figure

004780"5702960

7. The min term controller 74 then forwards the min term values (M\_STRU INFO) to the min term generator 76 and the equation core 78.

The min term generator 76 performs the actual min term comparisons between the min terms fetched by the min term controller and the selected byte of the incoming data stream. For example, the min term generator 76 simultaneously compares in Figure 7 the incoming data byte B1 with the min terms M1, M9, and M14, and provides the min term comparison results (mt\_result) to the equation core 78. During the next comparison cycle, the min term generator 76 simultaneously compares the incoming data byte B2 with the min terms M2, M10, and M15. According to the disclosed embodiment, the min term generator is configured for simultaneously comparing the incoming data stream to up to eight min terms.

The equation core 78 is configured for generating a frame tag based on the min term comparison results received from the min term generator 76, relative to the relevant templates 62. For example, the equation core 78 evaluates equation 1, illustrated in Figure 6, by evaluating the min term results sequentially as the results are supplied from the min term generator. For example, if the comparisons for each of the min terms M1, M2, M3, M4, M5, and M6 result in a true condition, described below with respect to Figure 8, then the end condition is matched in equation 1, causing the equation core 78 to generate a tag corresponding to the condition specified for equation 1. The frame tag identifies the nature of the incoming data packet, as well as the action that needs to be performed by the switch fabric 25.

Hence, as the packet arrives at the switch port 20, the packet classifier 24 immediately examines the arriving data at a programmable quantum (byte word, double word, etc.) boundary for a match against the appropriate field in the templates. If a match is found, the packet classifier 24 keeps track of which templates have matched and continues to examine the rest of the incoming packet fields in order of their arrival. Along the way, templates no longer matching the packet are discarded from further comparisons. At the end of the packet, or at the end of a fully matched template, the equation core 78 of the packet classifier 24 tags the packet with the tags that are associated with the matching templates. There may be circumstances when more than one template matches the packet. In such cases, a tag priority resolution device 81 performs a priority based determination, as defined by the user, so that at or near the end of frame, a single frame tag value is chosen and sent to the switch fabric 25.

Figure 8 is a diagram illustrating the data format of the min term structure in the min term memory 70. According to the disclosed embodiment, the min terms are stored in the min term memory 70 in an order corresponding to the relative position of the data byte to be compared. Hence,

09637015-081400



all min terms that are to be compared to the first data byte are stored together in a first part of the min term memory, followed by min terms to be compared with the second data byte, etc.

Alternatively, the min terms may be stored in an order based on relevant information within the IP header, as well as the relative position of the data byte to be compared. Hence, the min terms  
 5 may be stored in an order for comparing the sequence of data bytes providing the source IP address, destination IP address, and source and source and destination ports; in this case, non-relevant data bytes at the beginning of the IP frame would not have associated min terms stored at the beginning of the min term memory 70, further improving the efficiency of the min term memory 70.

Each table entry 90 includes a min term portion and an evaluation portion. The min term  
 10 portion includes a mask field (MASK) 94, an expected data field (EXP\_DATA) 96, and an operator field (OPERATOR) 98. Based on the position of the table entry 90 in the min term memory 70, the min term controller 74 is able to determine which byte of the IP packet 32 that needs to be compared with the corresponding min term, relative to the beginning 64 of the IP packet. The mask field 94 is a mask that is used by the min term generator 76 in performing comparisons; if the mask has a bit set to  
 15 1, the value is compared, and if the mask value has zeros in the field, the comparison is a don't care. The expected data field 96 specifies the expected data to be compared with the relevant data byte of the IP packet 32. The operator field 98 specifies the type of comparison to be performed by the min term generator, for example: less than, less than or equal to, equal to, greater than, greater than or equal to, and not equal to.

The evaluation portion includes a branches portion 100, a response portion (RINP1) 102 for  
 20 the case where the comparison of the min term portion is true, a second response portion (RINP0) 106 for the case where the comparison of the min term portion is false, and an equation identifier 110. The branches portion 100 specifies the order of the OR term in the equation; for example, the min term M1 as shown in Figures 6 and 7 would have its branches portion set to 0000 1111, indicating that  
 25 the first four branches of the equation specified in the equation identifier field 110 are to include the corresponding min term. The use of eight bits for the branches portion assumes that there are a maximum of eight branches in any given equation.

The response portion 102 specifies the operation to be performed if the min term portion is  
 30 evaluated as true relative to the compared data byte. In particular, the finish bit (FIN) is set to one if the results of the equation is determined if the min term result is true; the back to initial (BINIT) is set to one if the evaluation process should return to the initial state (init) if the min term result is true. For example, in the case of min term M1, the FIN bit and the BINIT bit of RINP1 are set to zero, since additional comparisons are needed if the min term result is true. In the case of min terms M5, M6,

09637015-081400

M7, and M8, the FIN bit of RINP1 is set to one, since a comparison result of “true” results in the end of the evaluation, as shown in Figure 6.

The response portion 106 specifies the operation to be performed if the min term portion is evaluated as false relative to the compared data byte. In particular, the finish bit (FIN) is set to one if the results of the equation is determined if the min term result is false; the back to initial (BINIT) is set to one if the evaluation process should return to the initial state (init) if the min term result is false. For example, in the case of min term M1, the FIN bit is set to zero and the BINIT bit of RINP1 is set to one, such that the equation would return to the INIT state if the min term result M1 was false, as shown in Figure 6.

The equation identifier field 110 identifies the equation (or template if there is only one template in an equation) that the min term corresponds to.

Hence, the equation core 78 determines whether any specified equation has a template 62 that matches the incoming data stream. Based on the multiple simultaneous comparisons of the incoming data stream with the multiple templates 62, the equation core 78 can identify a matching equation, and generate the appropriate tag corresponding to the matched equation for frame forwarding decisions at the switching fabric 25. If desired, the core 78 by also output a command to the header modifier 29 to modify the layer 2 header, the layer 3 header, or both, before transferring the data to the switch.

According to the disclosed embodiment, a network switch port includes a filter capable of performing multiple simultaneous comparisons between the incoming data stream of the data packet and multiple templates configured for identifying a corresponding protocol. Since the packet classifier module 24 can process any of the bytes of the IP frame 32, the packet classifier module 24 can interpret all the header information in the IP packet 32 from layer 3 up to layers 7 protocols. Moreover, the multiple simultaneous comparisons enables the network switch 12 to perform layer 3 switching for 100 Mbps and gigabit networks without blocking in the network switch. Finally, the multiple simultaneous comparisons in the order in which the data is received enables real time comparisons to be performed, as opposed to alternative schemes such as programmable logic arrays (PLAs), which would require the entire header to be received before processing can begin.

Hence, in accordance with the disclosed embodiment, the resolution of a packet classification and associated switch control attributes can occur as soon as the last bit of the packet arrives at the switch port, thereby reducing the overall switch latency. Thus, the packet classifier can be configured to forward packets based on any protocol, residing at any layer of the ISO model, at wire speed.

While this invention has been described with what is presently considered to be the most practical preferred embodiment, it is to be understood that the invention is not limited to the disclosed

09637015-081400

embodiments, but, on the contrary, is intended to cover various modifications and equivalent arrangements included within the spirit and scope of the appended claims.

004780" 5T07E960